

Cuestionario de teoría 2

Antonio Álvarez Caballero
analca3@correo.ugr.es

15 de mayo de 2016

1. Cuestiones

Cuestión 1. Sean x e y dos vectores de observaciones de tamaño N . Sea

$$\text{cov}(x, y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

la covarianza de dichos vectores, donde \bar{z} representa el valor medio de los elementos de z . Considere ahora una matriz X cuyas columnas representan vectores de observaciones. La matriz de covarianzas asociada a la matriz X es el conjunto de covarianzas definidas por cada dos de sus vectores columnas. Defina la expresión matricial que expresa la matriz $\text{cov}(X)$ en función de la matriz X .

Solución. Se define la matriz de covarianzas $(\text{cov}(X))_{ij}$ por componentes usando la definición de covarianza $\text{cov}(x, y)$ arriba indicada como:

$$(\text{cov}(X))_{ij} = \text{cov}(X_i, X_j)$$

con X_i la i -ésima columna de la matriz X .

Usando la igualdad $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$ de Statistical Inference de Casella y Berger (Teorema 5.2.4) es fácil ver:

$$\text{cov}(X_i, X_j) = \mathbb{E}[X_i X_j^T] - \mathbb{E}[X_i] \mathbb{E}[X_j]$$

Utilizando esta expresión, es directo ver que estas son las componentes de la matriz $\mathbb{E}[X X^T] - \mathbb{E}[X] \mathbb{E}[X^T]$, esto es:

$$(\mathbb{E}[X X^T] - \mathbb{E}[X] \mathbb{E}[X^T])_{ij} = \mathbb{E}[X_i X_j] - \mathbb{E}[X_i] \mathbb{E}[X_j]$$

Luego la expresión para la matriz de covarianzas es:

$$\text{cov}(X) = \mathbb{E}[X X^T] - \mathbb{E}[X] \mathbb{E}[X]$$

Cuestión 2. Considerar la matriz hat definida en regresión, $H = X(X^T X)^{-1} X^T$, donde X es una matriz $N \times (d+1)$, y $X^T X$ es invertible.

1. Mostrar que H es simétrica
2. Mostrar que $H^K = H$ para cualquier entero K

Solución. Resolvamos ambas partes:

1. Para ver que H es simétrica, no hay más que ver que es igual a su traspuesta. Hay que tener en cuenta que la traspuesta cambia el orden del producto matricial, además de que sabemos que $X^T X$ es simétrica.

$$H^T = \left(X(X^T X)^{-1} X^T \right)^T = X^{TT} \left(X(X^T X)^{-1} \right)^T = X(X^T X)^{-T} X^T = X(X^T X)^{-1} X^T = H$$

2. Para ver que $H^K = K \forall K \in \mathbb{N}$, lo haremos por inducción.

$K = 2$: Es claro que $H^2 = H$:

$$H^2 = X \underbrace{(X^T X)^{-1} X^T X}_{I} (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H$$

$K \xRightarrow{?} K + 1$: Suponiendo que es cierto para K :

$$H^{K+1} = H \Leftrightarrow \underbrace{H^K}_{\text{Inducción}} H = H \Leftrightarrow HH = H \Leftrightarrow \underbrace{H^2}_{K=2} = H$$

Lo cual ya hemos demostrado. Así ya lo hemos probado $\forall K \in \mathbb{N}$

Cuestión 3. Resolver el siguiente problema: Encontrar el punto (x_0, y_0) sobre la línea $ax + by + d = 0$ que este más cerca del punto (x_1, y_1) .

Solución. La función a minimizar es la función distancia, que en el plano \mathbb{R}^2 está definida por:

$$d(x, y) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Como la raíz cuadrada es una función creciente, podemos minimizar lo que queda dentro de la raíz, ya que alcanzan el mínimo en el mismo punto. Por tanto, nuestra función f a minimizar es:

$$f(x, y) = (x_1 - x_2)^2 + (y_1 - y_2)^2$$

Veamos ahora la restricción. El punto (x_0, y_0) que buscamos debe pertenecer a la recta dada luego la restricción g viene dada por:

$$g(x_0, y_0) = ax_0 + by_0 + d$$

Uniendo ambas expresiones, utilizando el método de los multiplicadores de *Lagrange*, nos queda:

$$\mathcal{L}(x_0, y_0, \lambda) = (x_1 - x_0)^2 + (y_1 - y_0)^2 - \lambda(ax_0 + by_0 + d)$$

Ahora debemos calcular el gradiente de este campo escalar, igualar al vector $\vec{0}$ y resolver el sistema de ecuaciones resultante.

$$\frac{\partial \mathcal{L}}{\partial x_0} = 0 \Leftrightarrow -2(x_1 - x_0) + \lambda a = 0 \Leftrightarrow x_0 = x_1 - a \frac{\lambda}{2}$$

$$\frac{\partial \mathcal{L}}{\partial y_0} = 0 \Leftrightarrow -2(y_1 - y_0) + \lambda b = 0 \Leftrightarrow y_0 = y_1 - b \frac{\lambda}{2}$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = 0 \Leftrightarrow ax_0 + by_0 + d = 0 \Leftrightarrow a \left(x_1 - \frac{\lambda a}{2} \right) + b \left(y_1 - \frac{\lambda b}{2} \right) + d = 0 \Leftrightarrow \lambda = \frac{2ax_1 + 2by_1 + 2d}{a^2 + b^2}$$

Ahora sustituimos λ en las equivalencias que hemos obtenido para x_0 e y_0 :

$$x_0 = x_1 - a \frac{ax_1 + by_1 + d}{a^2 + b^2}$$

$$y_0 = y_1 - b \frac{ax_1 + by_1 + d}{a^2 + b^2}$$

Por tanto, el punto (x_0, y_0) de la recta $ax + by + d = 0$ más cercano a un punto del plano (x_1, y_1) es

$$(x_0, y_0) = \left(x_1 - a \frac{ax_1 + by_1 + d}{a^2 + b^2}, y_1 - b \frac{ax_1 + by_1 + d}{a^2 + b^2}\right)$$

Cuestión 4. Consideremos el problema de optimización lineal con restricciones definido por

$$\begin{aligned} & \text{Min}_z \mathbf{c}^T \mathbf{z} \\ & \text{Sujeto a } \mathbf{A} \mathbf{z} \leq \mathbf{b} \end{aligned}$$

donde \mathbf{c} y \mathbf{b} son vectores y \mathbf{A} es una matriz.

1. Para un conjunto de datos linealmente separable mostrar que para algún w se debe de verificar la condición $\mathbf{y}_n \mathbf{w}^T \mathbf{x}_n > 0$ para todo $(\mathbf{x}_n, \mathbf{y}_n)$ del conjunto.
2. Formular un problema de programación lineal que resuelva el problema de la búsqueda del hiperplano separador. Es decir, identifique quienes son \mathbf{A} , \mathbf{z} , \mathbf{b} y \mathbf{c} para este caso.

Solución. Resolvemos por partes:

1. La condición $\mathbf{y}_n \mathbf{w}^T \mathbf{x}_n > 0$ quiere decir que el dato (x_n, y_n) está bien clasificado, como ya sabemos para clasificación binaria. Entonces, demostramos por contrarrecíproco ($(a \Rightarrow b) \Leftrightarrow (\neg b \Rightarrow \neg a)$): Si no existe w que cumpla $\forall (\mathbf{x}_n, \mathbf{y}_n) \mathbf{y}_n \mathbf{w}^T \mathbf{x}_n > 0$, eso quiere decir que no clasifica bien todos los datos, luego el conjunto no es linealmente separable.

Cuestión 5. Probar que en el caso general de funciones con ruido se verifica que $\mathbb{E}_{\mathcal{D}}[E_{out}] = \sigma^2 + \mathbf{bias} + \mathbf{var}$ (ver transparencias de clase).

Solución. Nuestro problema se reduce a calcular, tal como se hizo en clase, la expresión $\mathbb{E}_{\mathcal{D}}[E_{out}(g^D(x))]$, donde la función objetivo f tiene ruido. Suponemos, como siempre, que el ruido ϵ tiene media 0 y varianza σ^2 .

$$\mathbb{E}_{\mathcal{D}} [E_{out}(g^D(x))] = \mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_x \left[\left(g^D(x) - y(x) \right)^2 \right] \right] = \mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_x \left[\left(g^D(x) - f(x) - \epsilon(x) \right)^2 \right] \right]$$

$$\mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_x \left[\left(g^D(x) - f(x) - \epsilon(x) \right)^2 \right] \right] = \mathbb{E}_x \left[\mathbb{E}_{\mathcal{D}} \left[\left(g^D(x) - f(x) - \epsilon(x) \right)^2 \right] \right]$$

Desarrollando el trinomio queda:

$$\mathbb{E}_x \left[\mathbb{E}_{\mathcal{D}} \left[g^D(x)^2 + f(x)^2 + \epsilon(x)^2 - 2f(x)g^D(x) - 2g^D(x)\epsilon(x) - 2f(x)\epsilon(x) \right] \right]$$

Ahora sumamos y restamos $\bar{g}(x)^2$ para obtener las expresiones para $\mathbf{bias}(x)$ y $\mathbf{var}(x)$.

$$\mathbb{E}_x \left[\mathbb{E}_{\mathcal{D}} \left[g^D(x)^2 - \bar{g}(x)^2 + \bar{g}(x)^2 + f(x)^2 + \epsilon(x)^2 - 2f(x)g^D(x) - 2g^D(x)\epsilon(x) - 2f(x)\epsilon(x) \right] \right]$$

Usamos la linealidad de la esperanza y reordenamos para obtener dichas expresiones. La varianza es la indicada porque el doble producto de que debe salir del binomio al cuadrado es 0, ya que $\bar{g}(x)$ no depende de \mathcal{D} .

$$\mathbb{E}_x \left[\underbrace{\mathbb{E}_{\mathcal{D}} [g^D(x)^2] - \bar{g}(x)^2}_{var(x)} + \underbrace{\bar{g}(x)^2 - 2f(x) \mathbb{E}_{\mathcal{D}} [g^D(x)] + f(x)^2}_{bias(x)} - 2\epsilon(x) \underbrace{\mathbb{E}_{\mathcal{D}} [g^D(x)]}_{\bar{g}(x)} - 2f(x)\epsilon(x) + \epsilon(x)^2 \right]$$

Entonces nos queda

$$\mathbb{E}_{\mathcal{D}}[E_{out}(g^D(x))] = \mathbb{E}_x [var(x) + bias(x) - 2\epsilon(x)\bar{g}(x) - 2f(x)\epsilon(x) + \epsilon(x)^2]$$

Ahora usamos de nuevo la linealidad de la esperanza. Además, destacar que $f(x)$ y $\epsilon(x)$ son independientes, luego la esperanza del producto es el producto de las esperanzas.

$$\mathbb{E}_{\mathcal{D}}[E_{out}(g^D(x))] = var + bias - 2 \underbrace{\mathbb{E}_x [\epsilon(x)] \mathbb{E}_x [\bar{g}(x)]}_0 - 2 \mathbb{E}_x [f(x)] \underbrace{\mathbb{E}_x [\epsilon(x)]}_0 + \underbrace{\mathbb{E}_x [\epsilon(x)^2]}_{\sigma^2}$$

Con esto, ya lo tenemos: por hipótesis, sabemos que la media del ruido es 0 y la varianza σ^2 , luego nos queda la expresión que buscábamos:

$$\mathbb{E}_{\mathcal{D}}[E_{out}(g^D(x))] = var + bias + \sigma^2$$

Cuestión 6. Consideremos las mismas condiciones generales del enunciado del Ejercicio.2 del apartado de Regresión de la relación de ejercicios.2. Considerar ahora $\sigma = 0,1$ y $d = 8$, ¿cual es el más pequeño tamaño muestral que resultará en un valor esperado de E_{in} mayor de 0,008?.

Solución. Conociendo la expresión del valor esperado del E_{in} y acotándolo inferiormente por 0,008:

$$\mathbb{E}_{\mathcal{D}} [E_{in}(w_{lim})] = \sigma^2 \left(1 - \frac{d+1}{N} \right) \geq 0,008$$

Ahora sólo sustituimos los valores que nos han dado y despejamos el tamaño de la muestra N .

$$0,1^2 \left(1 - \frac{9}{N} \right) \geq 0,008 \Rightarrow \left(1 - \frac{9}{N} \right) \geq 0,8 \Rightarrow -\frac{9}{N} \geq -0,2 \Rightarrow \frac{9}{0,2} \leq N \Rightarrow 45 \leq N$$

Así, el mínimo N para que el valor esperado de E_{in} sea 0,008 es 45.

Cuestión 7. En regresión logística mostrar que

$$\nabla E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{y_n w^T x_n}} = \frac{1}{N} \sum_{n=1}^N -y_n x_n \sigma(-y_n w^T x_n)$$

Argumentar que un ejemplo mal clasificado contribuye al gradiente más que un ejemplo bien clasificado.

Solución. Partimos de la expresión de E_{in} conocida:

$$E_{in}(w) = \frac{1}{N} \sum_{n=0}^N \ln(1 + e^{-y_n w^T x_n})$$

Ahora sólo calculamos su gradiente con respecto a w , por lo que coincide con su parcial con w .

$$\nabla E_{in}(w) = \frac{\partial}{\partial w} \left(\frac{1}{N} \sum_{n=0}^N \ln(1 + e^{-y_n w^T x_n}) \right) = \frac{1}{N} \sum_{n=1}^N \frac{-y_n x_n e^{-y_n w^T x_n}}{1 + e^{-y_n w^T x_n}} = \frac{1}{N} \sum_{n=1}^N \frac{-y_n x_n}{1 + e^{y_n w^T x_n}}$$

$$\nabla E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{y_n w^T x_n}} = \frac{1}{N} \sum_{n=1}^N -y_n x_n \sigma(-y_n w^T x_n)$$

Es claro que una solución mal clasificada influye más que una buena. Una solución bien clasificada equivale a $y_n w^T x_n > 0$, luego la exponencial del denominador será mayor que 1. Al dividir por algo más grande, el cociente se hace más pequeño. En el caso de $y_n w^T x_n < 0$, la exponencial será una exponencial negativa, con valores menores a 1. Luego el denominador será más pequeño, y el cociente más grande. Luego la sumatoria será más grande.

Cuestión 8. Definamos el error en un punto (x_n, y_n) por

$$e_n(w) = \max(0, -y_n w^T x_n)$$

Argumentar que el algoritmo PLA puede interpretarse como SGD sobre e_n con tasa de aprendizaje $\eta = 1$.

Solución. Recordamos que la función de actualización de los pesos para el SGD es

$$w = w - \eta \cdot \nabla e_n(w)$$

Además, la función de error está bien definida. Si un dato está mal clasificado, $y_i w^T x_i > 0$, por lo cual el máximo de 0 y $-y_i w^T x_i$ es 0. En caso de estar mal clasificado, el máximo es dicho valor.

Vamos a calcular el gradiente del error y veamos qué sale.

$$\nabla e_n(w) = \frac{\partial}{\partial w} \left(\max(0, -y_n w^T x_n) \right) = \max(0, -y_n x_n)$$

Entonces, para datos mal clasificados, la función de actualización de los pesos es:

$$w = w - \eta \cdot (-y_n x_n) \xrightarrow{\eta=1} w = w + 1 \cdot y_n x_n$$

Que es justo la fórmula para actualizar pesos en el PLA.

Cuestión 9. El ruido determinista depende de \mathcal{H} , ya que algunos modelos aproximan mejor f que otros.

1. Suponer que \mathcal{H} es fija y que incrementamos la complejidad de f .
2. Suponer que f es fija y decrementamos la complejidad de \mathcal{H}

Contestar para ambos escenarios: ¿En general subirá o bajará el ruido determinista? ¿La tendencia a sobreajustar será mayor o menor? (Ayuda: analizar los detalles que influyen el sobreajuste)

Solución. Analizamos por partes:

1. El ruido determinista está directamente determinado por \mathcal{H} , ya que es la parte que no ajustamos bien de f . Si incrementamos la complejidad de f , el ruido determinista aumentará, ya que al aumentar la complejidad de la función objetivo manteniendo nuestro espacio de hipótesis, cada vez ajustaremos peor. Luego la diferencia entre la función objetivo (cuya complejidad va aumentando) y nuestra función estimada, que pertenece a \mathcal{H} (que dejamos fijo) será cada vez mayor.

Acerca del sobreajuste, la tendencia será a sobreajustar más. Esto se debe a que se intentará aprender todo lo posible de f , sin poder determinar si algunos datos son ruido o no.

2. En este caso también va a subir el ruido determinista. Porque la f es fija, pero al quitar complejidad a \mathcal{H} cada vez podremos ajustar peor cada uno de los datos de la muestra, teniendo un ajuste peor de f y por tanto más ruido determinista.

Acerca del sobreajuste, la tendencia será a sobreajustar menos. Al decrementar la complejidad de \mathcal{H} estamos ajustando peor cada uno de los puntos de la muestra, como ya hemos comentado. Al no poder focalizar sobre cada uno de los puntos individualmente, el sobreajuste baja, aunque cada vez ajustará peor la función objetivo.

Cuestión 10. La técnica de regularización de Tikhonov es bastante general al usar la condición

$$w^T \Gamma^T \Gamma w \leq C$$

que define relaciones entre las w_i (La matriz Γ_i se denomina regularizador de Tikhonov)

1. Calcular Γ cuando $\sum_{q=0}^Q w_q^2 \leq C$
2. Calcular Γ cuando $(\sum_{q=0}^Q w_q)^2 \leq C$

Argumentar si el estudio de los regularizadores de Tikhonov puede hacerse a través de las propiedades algebraicas de las matrices Γ .

Solución. Veamos ambos casos:

1. En este caso Γ debe ser cualquier matriz ortogonal, es decir, que cumpla que su inversa coincide con su traspuesta.
2. En este caso es un poco más complicado. Vamos a intentar sacarlo para un caso pequeño, $Q = 3$. Para este caso, debemos buscar una matriz Γ que haga cumplir la condición

$$(w_1 + w_2 + w_3)^2 \leq C$$

Por tanto, debemos buscar Γ que haga que $w^T \Gamma^T = w_1 + w_2 + w_3 = \Gamma w$ en alguna de sus componentes (necesariamente la misma, para que al realizar el producto escalar consigo misma dé dicho resultado al cuadrado). Sabemos que w es columna y por tanto w^T fila. Entonces es claro que

$$\Gamma w = (w_1 + w_2 + w_3, 0, 0) \Leftrightarrow \Gamma = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Para el otro caso, $w^T \Gamma^T$, veamos si pasa lo mismo.

$$w^T \Gamma^T = (w_1, w_2, w_3) \cdot \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} = (w_1 + w_2 + w_3, 0, 0)$$

Entonces nos queda que

$$w^T \Gamma^T \Gamma w = (w_1 + w_2 + w_3, 0, 0) \cdot \begin{pmatrix} w_1 + w_2 + w_3 \\ 0 \\ 0 \end{pmatrix} = (w_1 + w_2 + w_3)^2$$

Que es justo lo que queríamos. Para el caso general, es exactamente igual, tomando vectores y matrices cuadradas de dimensión Q en vez de 3.

BONUS. Considerar la matriz $\hat{H} = X(X^T X)^{-1} X^T$. Sea X una matriz $N \times (d+1)$, y $X^T X$ invertible. Mostrar que $\text{traza}(\hat{H}) = d+1$, donde traza significa la suma de los elementos de la diagonal principal. (+1 punto)

Solución. Para resolver esto sólo necesitamos utilizar una propiedad básica de álgebra lineal:

$$\text{Traza}(AB) = \text{Traza}(BA)$$

Utilizando esto, dividimos \hat{H} en dos matrices A y B .

$$A = X(X^T X)^{-1}$$

$$B = X^T$$

Es claro que $AB = \hat{H}$. Es claro que BA es la matriz identidad. Lo único que hay que tener cuidado es con las dimensiones de esta matriz.

$$X \in \mathcal{M}_{N \times d+1} \Rightarrow X^T X \in \mathcal{M}_{d+1 \times d+1}$$

Como una matriz y su inversa tienen la misma dimensión, entonces:

$$BA = \underbrace{X^T X}_M \underbrace{(X^T X)^{-1}}_{M^{-1}} = I_{d+1}$$

Y la traza de esta matriz es $d+1$. Como coincide con la traza de \hat{H} , podemos concluir que

$$\text{traza}(\hat{H}) = d+1$$